# The Higgs Boson Prediction CS-433 Machine Learning – Project 1

Antoine Clivaz, Trung-Dung Hoang, Lorenzo Taschin School of Computer and Communication Sciences, EPFL, Switzerland

Abstract—A critical problem encountered when studying the Higgs Boson is to distinguish between the signal of Higgs Boson and other processes by utilizing incomplete features. In this report, we use Machine Learning to resolve the above problem by considering several models and data processing approaches. To this end, we discover that dividing the dataset into multiple subsets which have the same available features and then applying further feature augmentation yields the best performance, which was 81.7% accuracy and 0.725 F1-score on Alcrowd.com.

# I. INTRODUCTION

In the Standard Model of particle physics, the Higgs Boson is a fundamental element that has the role of giving mass to other elementary particles.

A Higgs Boson can be generated as a by-product of a proton-proton collision. However, it decays rapidly into other particles, so scientists have to measure decay signatures to determine whether it was the result of a Higgs Boson or some other process. The problem becomes more challenging as many decay signatures look similar.

In this project, our goal is to tackle the problem with a Machine Learning approach. We concentrate on examining multiple ways to handle missing signatures, then compare multiple models to find the best method for Higgs Boson predictions.

## II. DATA PREPROCESSING

#### A. Data observations

In this project, the training dataset contains 250.000 pairs of features and binary labels. For each sample, there are a total of 30 features that can be divided into two categories<sup>1</sup>: The PRI (primitives) category that represents quantities measured by the detector and the DER (derived) category that are quantities computed from the primitives features.

- 1) Missing values: The occurrences of the value -999, which appears in several features, corresponding to the values that are meaningless or cannot be computed. These missing values can be categorized into two groups:
  - **Missing at random (MAR)**: The propensity for a feature to be missing depends on the observed data, in this case, *PRI\_jet\_num*. There are multiple features whose definitions depend on the number of jets.

- Missing not at random (MNAR): The value of a feature that is missing depends on the reason why it is missing. In this case, the feature *DER\_mass\_MMC* may be undefined if the topology of the event is too far from the expected topology<sup>1</sup>.
- 2) Outliers: Outliers can negatively impact our linear model so it is important to detect them if possible. By simply visualizing the distribution of each feature, excluding missing values, we can observe how heavy-tailed some of them are, which indicates the presence of outliers in the data.

#### B. Data Preparation

According to observations we made in the above section, we decide to examine several approaches to preprocess the data.

- Handling missing values There are 11 out of 30 features containing missing values. We consider two strategies. The first strategy is removing features that contain missing values. The hypothesis is that the remaining features are still relevant enough to predict the outcome. The second strategy consists in replacing missing values with the mean/median/most-frequent value of these features.
- Splitting the data We also decide to split the data into 4 distinct groups according to PRI\_jet\_num value. For each group, we get rid of features containing only one value, including PRI\_jet\_num and other features that have only missing values. Finally, we handle missing values for the DER\_mass\_MMC feature by the second approach mentioned above.
- Furthermore, we can extend the idea by splitting the data into 8 distinct groups, that are completely free of missing values, according to *PRI\_jet\_num* value and whether *DER\_mass\_MMC* is defined or not. Then, we drop features whose values are the same for every sample.
- Features normalization and outliers removal The data is centered on zero mean and unit variance to avoid the ill-condition problem. Transformed values can be considered as z\_score and can be marked as outliers if they are larger than 3. Any sample containing at least one outlier feature value is dropped.
- Features augmentation We then apply polynomial expansion of degree up to k ( $k \ge 2$ ) on the features

<sup>&</sup>lt;sup>1</sup>https://higgsml.ijclab.in2p3.fr/files/2014/04/documentation\_v1.8.pdf

of each dataset. More precisely, the original features  $\{x_i\}_{1\leq i\leq d}$  are transformed to  $\{x_i^j\}_{1\leq i\leq d, 2\leq j\leq k}$ . Furthermore, we also consider creating interactions between features by the following mapping:  $\{x_i\}_{1\leq i\leq d}\mapsto \{\prod_{x_j\in A}x_j|A\in {x_i\}\choose k}\}$  where  ${x_i\choose k}$  is the set all k-combinations of  $\{x_i\}$ . All these augmented features are concatenated to the original ones to form a whole.

#### III. MODELS AND METHODS

We utilize 4 out of 6 different training algorithms obtained from the lab: least squares (**LS**), ridge regression (**RR**), linear regression (**LR**) and regularized logistic regression (**RLR**). For simplicity, we keep 12 regularization fixed and apply grid search only for the weight of regularizer (*lambda*) of two regularized models.

To avoid combinatorial explosion, we use the same stepsize scheduler for every algorithm rather than tuning this hyperparameter. More specifically, while training, we keep track of the smallest loss and the number of steps the algorithm takes without reaching a loss being less than the current smallest loss by an amount of tol. If this number is larger than a predefined value, the step size seems to be too large, so we decrease it r times.

The algorithm terminates if it reaches the maximum number of steps or the step size is smaller than a predefined threshold.

## IV. EXPERIMENT

#### A. Experiment setting

For each setting, we use 5-fold to evaluate the model. We report the average accuracy and fl\_score for 5 runs and use them to compare model performance. In the case of data splitting, we train different models for each split.

The initial value of step size, the maximum number of steps without progress, the *tol* value, the r value, the smallest acceptable step size are 1.0, 3, 1e-5, 5 and 1e-3, respectively. The maximum of training steps is 1000. For *gamma*, we search the best value among 10 values spaced evenly on a log scale range from 1e-4 to 1.

The features are always added bias terms and normalized to have zero mean and unit variance. The weights of the models are initialized by vectors of ones.

#### B. Experiment results

- 1) Missing value handler: In this part, we want to investigate the effect of missing value on performance. We compare 5 settings:
  - (M1) Drop features containing missing values
  - (M2) Replacing with mean
  - (M3) Spliting by mass
  - (M4) Spliting by jet\_num
  - (M5) Spliting by jet\_num and mass

After handling missing data, features are normalized and added bias terms as mentioned in the setting session. Our hypothesis is that all the remaining features are meaningful, removing or changing them can make us lose useful information and can lead to poor performance. The results shown in table I confirm this hypothesis. The more you modify the original values, the worse the model is. Furthermore, we observe that logistic regression performs slightly better than linear regression.

Table I
MISSING VALUE HANDLER

	LS	RR	LR	RLR
M1	73.342±0.264	$73.338 \pm 0.266$	73.857±0.247	$73.872 \pm 0.264$
<b>M2</b>	$74.422 \!\pm 0.210$	$74.438 \pm 0.208$	$0.750 \pm 0.259$	$75.023 \pm 0.258$
M3	$75.001\pm0.199$	$75.014 \pm 0.214$	$75.356 \pm 0.209$	$75.418 \!\pm 0.198$
<b>M4</b>	$75.872 \pm 0.181$	$75.875 \pm 0.170$	$76.460 \pm 0.142$	$76.468 \pm 0.132$
M5	$76.506 \pm 0.125$	$76.498 \pm 0.128$	$76.854 \!\pm 0.127$	$76.852 \pm 0.160$

2) Further features preprocessing: From now, due to the space limitation, we only demonstrate the results of logistic regression and its variant with (M5). We study the effect of outliers removal (OR) and polynomial expansion (PE) on the performance. In this experiment, we fix the upper bound of z\_score of normal values to 3. In terms of the degree of polynomial expansion, we set the degree to 3 as the result of grid search processes between 2 to 5. Table II shows that while removing outliers slightly helps improve performances, adding augmented features can significantly boost the accuracy and a combination of them yields the best result.

Table II FURTHER FEATURES PREPROCESSING

OR	PE	OR+PE
 	$80.940 \pm 0.121$ $81.354 \pm 0.109$	

In addition, we attempt to use interactions between features as well as more complex non-linear transformations like logarithm, sin, but the performance drops considerably. We hypothesize that the features themselves already contain interactions and non-linear transformations between primitive quantities, combining only primitive quantities may enhance the performance but utilizing every possible combination can be ineffective.

## V. CONCLUSION

In this project, we evaluate multiple learning models while paying attention to data preprocessing. We notice that handling data properly can exert positive impacts on the performance of the model. After several experiments, we obtain that a logistic regression model with 8 different subsets combined with outliers removal, feature polynomial expansion and normalization produces the best results:81.7% accuracy and 0.725 F1-score on Alcrowd.com.